

UTILITY PATENT APPLICATION TRANSMITTAL (Large Entity)

(Only for new nonprovisional applications under 37 CFR 1.53(b))

Docket No.
ARC9-2000-0079USI

Total Pages in this Submission

TO THE ASSISTANT COMMISSIONER FOR PATENTS

Box Patent Application
Washington, D.C. 20231

Transmitted herewith for filing under 35 U.S.C. 111(a) and 37 C.F.R. 1.53(b) is a new utility patent application for an invention entitled:

A METHOD FOR ADAPTING A K-MEANS TEXT CLUSTERING TO EMERGING DATA

and invented by:

William Scott Spangler

If a **CONTINUATION APPLICATION**, check appropriate box and supply the requisite information:

☐ Continuation ☐ Divisional ☐ Continuation-in-part (CIP) of prior application No.: _____

Which is a:

☐ Continuation ☐ Divisional ☐ Continuation-in-part (CIP) of prior application No.: _____

Which is a:

☐ Continuation ☐ Divisional ☐ Continuation-in-part (CIP) of prior application No.: _____

Enclosed are:

Application Elements

1. ☒ Filing fee as calculated and transmitted as described below
2. ☒ Specification having 30 pages and including the following:
 - a. ☒ Descriptive Title of the Invention
 - b. ☐ Cross References to Related Applications (if applicable)
 - c. ☐ Statement Regarding Federally-sponsored Research/Development (if applicable)
 - d. ☐ Reference to Microfiche Appendix (if applicable)
 - e. ☒ Background of the Invention
 - f. ☒ Brief Summary of the Invention
 - g. ☒ Brief Description of the Drawings (if drawings filed)
 - h. ☒ Detailed Description
 - i. ☒ Claim(s) as Classified Below
 - j. ☒ Abstract of the Disclosure

UTILITY PATENT APPLICATION TRANSMITTAL (Large Entity)

(Only for new nonprovisional applications under 37 CFR 1.53(b))

Docket No.
ARC9-2000-0079USI

Total Pages in this Submission

Application Elements (Continued)

3. ☒ Drawing(s) (when necessary as prescribed by 35 USC 113)
- a. ☒ Formal Number of Sheets 3 (Figures 1-3)
- b. ☐ Informal Number of Sheets _____
4. ☒ Oath or Declaration
- a. ☒ Newly executed (original or copy) ☐ Unexecuted
- b. ☐ Copy from a prior application (37 CFR 1.63(d)) (for continuation/divisional application only)
- c. ☒ With Power of Attorney ☐ Without Power of Attorney
- d. ☐ DELETION OF INVENTOR(S)
Signed statement attached deleting inventor(s) named in the prior application,
see 37 C.F.R. 1.63(d)(2) and 1.33(b).
5. ☐ Incorporation By Reference (usable if Box 4b is checked)
The entire disclosure of the prior application, from which a copy of the oath or declaration is supplied under
Box 4b, is considered as being part of the disclosure of the accompanying application and is hereby
incorporated by reference therein.
6. ☐ Computer Program in Microfiche (Appendix)
7. ☐ Nucleotide and/or Amino Acid Sequence Submission (if applicable, all must be included)
- a. ☐ Paper Copy
- b. ☐ Computer Readable Copy (identical to computer copy)
- c. ☐ Statement Verifying Identical Paper and Computer Readable Copy

Accompanying Application Parts

8. ☒ Assignment Papers (cover sheet & document(s))
9. ☐ 37 CFR 3.73(B) Statement (when there is an assignee)
10. ☐ English Translation Document (if applicable)
11. ☒ Information Disclosure Statement/PTO-1449 ☒ Copies of IDS Citations
12. ☐ Preliminary Amendment
13. ☒ Acknowledgment postcard
14. ☐ Certificate of Mailing
- ☐ First Class ☐ Express Mail (Specify Label No.): _____

A METHOD FOR ADAPTING A K-MEANS TEXT CLUSTERING TO

EMERGING DATA

BACKGROUND OF THE INVENTION

5 *Field of the Invention*

The present invention generally relates to computerized datasets and more particularly to a method and system for automatically categorizing, indexing, and classifying items within such datasets.

10 *Description of the Related Art*

K-means is a well known algorithm for clustering (i.e. partitioning) a dataset of numeric vectors, where each numeric vector has dimensionality M and there are N such vectors. The value, K , refers to an input parameter of the algorithm that determines the number of such clusters (i.e. partitions) that the algorithm will produce at completion. In general, K-means, from a given starting point, finds a locally optimum way to cluster the dataset into K partitions so as to minimize the average difference between the mean of each cluster (cluster centroid) and every member of that cluster. This difference is measured by some distance metric, such as Euclidean distance.

20

In the case of text datasets, the N vectors represent text documents and dimensionality M refers to the occurrence of certain keywords or phrases in the text documents. The dictionary of keywords or phrases may be derived by counting the occurrence of all words and/or phrases in the text corpus and
5 selecting those words and phrases that occur most often.

Problems to be Solved by the Invention

The major drawback of conventional techniques is that they cannot take
10 an existing K-means classification (as represented by the K centroids) and adapt it to a new, but related, dataset (e.g. one taken from the same domain at a later time).

A typical instance of this problem arises with helpdesk problem tickets. Support for consumer and commercial products is often provided telephonically.
15 In such situations, an operator (often referred to as a "helpdesk" operator) receives a telephone call with a problem. The telephone operator assigns each problem a specific identification code and records the user's problem (and the help advice provided) in a computerized file. In such a system, if the user calls back, other help desk operators can retrieve the computerized file using the specific
20 identification code. This prevents the user from having to wait to speak with a

specific help desk operator. Additionally, this system provides a helpdesk dataset of problems and solutions which other help desk operators can access when offering potential solutions to users.

Free form computer helpdesk datasets consist primarily of short text descriptions, composed by the helpdesk operator for the purpose of summarizing what problem a user had and what was done by the helpdesk operator to solve that problem. A typical text document (known as a problem ticket) from this data consists of a series of exchanges between an end user and an expert helpdesk advisor.

For example, one problem ticket may read as follows: "1836853 User calling in with WORD BASIC error when opening files in word. Had user delete NORMAL.DOT and had the user reenter Word. The user was fine at that point. 00:04:17 ducar May 2:07:05:656P". This problem ticket begins with the unique identification number, which is followed by a brief identification of the user's problem, the solution offered, the help desk operators name or identification symbol, and a date and time stamp.

Problem tickets may include only of a single symptom and resolution pair, as in the above example, or the problem tickets may span multiple questions, symptoms, answers, attempted fixes, and resolutions, all pertaining to the same basic issue. Problem tickets are opened when the user makes the first call to the helpdesk and are closed when all user problems documented in the

first call are finally resolved in some way. Helpdesk operators enter problem tickets directly into the database. Spelling, grammar and punctuation are inconsistent. The style is terse and the vocabulary very specialized. Therefore, the contents of the help desk dataset are not very useful for answering future

5 problems because the data is not easily searched, categorized, or indexed.

In co-pending U.S. Patent Application 09/542,859, incorporated herein by reference, one solution to the foregoing problem involves mining helpdesk datasets for "problem resolution nuggets". In other words, this solution finds those problem tickets that describe both the symptoms of a problem and how the

10 helpdesk operator solved that problem. Such problem resolution nuggets can then be utilized to create solution documents to put on a web site, or to create automated knowledge bases that can diagnose the user's problem automatically and suggest a corrective action.

With such a system, helpdesk problem tickets may be classified using a

15 K-means algorithm for the month of January. Subsequently, there may be a need to reuse this-classification in the month of February. It is assumed that the underlying classes in the dataset have not changed drastically between January and February, but certain small changes could have taken place. For example, new kinds of problems may have been introduced that necessitate new classes

20 being created to represent them.

However, conventional solutions merely run the K-means again from scratch on the new February data. Unfortunately, due to the fact that the resulting K-means clusters are highly dependent on the initial starting point (seeds), a strong possibility exists that the new February clustering will not have a very close relationship to the original K-clusters. In other words, for any given cluster in the original K-means clustering (January), there will not necessarily be a similar cluster of dataset points in the subsequent K-means clustering (February). This lack of continuity is a drawback when one is interested in tracking changes in the data (trends) over time.

To solve this problem, the invention identifies new emerging concepts in the succeeding dataset while retaining those concepts that carry over from the previous dataset.

SUMMARY OF THE INVENTION

It is, therefore, an object of the present invention to provide a structure and method of clustering documents in datasets which include clustering first documents and a first dataset to produce first document classes, creating centroid seeds based on the first document classes, and clustering second documents in a second dataset using the centroid seeds, wherein the first dataset and the second dataset are related.

The clustering of the first documents in the first dataset forms a first dictionary of most common words in the first dataset and generates a first vector space model by counting, for each word in the first dictionary, a number of the first documents in which the word occurs, and clusters the first documents in the first dataset based on the first vector space model, and further generates a second vector space model by counting, for each word in the first dictionary, a number of the second documents in which the word occurs.

Creation of the centroid seeds includes classifying second vector space model using the first document classes to produce a classified second vector space model and determining a mean of vectors in each class in the classified second vector space model, the mean includes the centroid seeds.

The invention also includes a method of clustering documents in datasets which form a second dictionary of most common words in the second dataset, generating a third vector space model by counting, for each word in the second dictionary, a number of the second documents in which the word occurs, and clustering the documents in the second dataset based on the second vector space model to produce a second dataset cluster, wherein the clustering of the second documents in the second dataset using the centroid seeds produces an adapted dataset cluster and, the method further including comparing classes in the adapted dataset cluster to classes in the second dataset cluster, and adding classes to the adapted dataset cluster based on the comparing.

003360-0369960

The invention can also include a system for clustering documents in datasets including a storage having a first dataset and a second dataset, a cluster generator operative to cluster first documents in the first dataset and produce first document classes, and a centroid seed generator operative to generate centroid seeds based on the first document classes, wherein the cluster generator clusters second documents in the second dataset using the centroid seeds.

ADVANTAGES OF THE INVENTION

One advantage of the invention lies in the ability to identify new emerging concepts in a succeeding dataset while retaining those concepts that carry over from a previous dataset.

BRIEF DESCRIPTION OF THE DRAWINGS

The foregoing and other objects, aspects and advantages will be better understood from the following detailed description of a preferred embodiment of the invention with reference to the drawings, in which:

Figure 1 is a flowchart illustrating one embodiment of the invention;

Figure 2 is a schematic architectural diagram of one embodiment of the invention; and

Figure 3 is a hardware embodiment for implementing the invention.

DETAILED DESCRIPTION OF PREFERRED EMBODIMENTS OF THE INVENTION

In the following description it is assumed that an initial text dataset, T1 (e.g., January helpdesk data), is classified first, followed by a new, but related text
5 dataset, T2 (e.g., February helpdesk data) which should be classified similarly, but should also be indexed to highlight emerging trends.

Referring now to the drawings, Figure 1 is a flowchart illustrating a first aspect of the invention and Figure 2 is a schematic diagram of the functioning elements of this embodiment of the invention. In item 100, the invention begins
10 by generating a first dictionary 206, D1, of frequently used words from dataset T1 200 using a dictionary generator 204. The most frequently occurring words in the corpus make up the dictionary. This reduced set of words will be used to compose a simple description of each document in the corpus. The number of words to be included in the dictionary is a user specified parameter.

15 Then, the vector space model generator 210 counts, for each word in the first dictionary D1 206, the number of documents in which the word in question appears, to produce a T1-D1 vector space mode 212. In item 102, by counting the occurrences of dictionary words in documents of dataset T1, the invention generates a matrix of non-negative integers where each column corresponds to a
20 word in the dictionary and each row corresponds to an example in the text corpus.

by the classifier 218 (based on the K-classes from the T1-D1 cluster) and inputs the centroid seeds to the K-means cluster generator 222, as shown in item 114.

The initial centroid (seed) for each class is found by summing up the columns of all examples in the class and dividing these values by the number of elements in the class. Centroid seeds are used to generate the initial clustering which is then optimized using the K-means approach.

More specifically, given a set of K centroids the invention finds, for each example point, the centroid to which it is closest. Each example point is now said to "belong to the class" of its nearest centroid. Then, the invention calculates the mean of each class. Each mean now becomes the new centroid, and there are again K centroids. The invention repeat the foregoing until no change results in class membership as the result of finding new centroids.

Obviously the starting place of this process will, to some extent, determine the final resulting classes. The starting points (seeds, or initial centroids) are conventionally selected in some random fashion to prevent undo biasing of the algorithm. To the contrary, the invention intentionally biases the algorithm towards the previous classification centroids. Thus, the invention directs the K-means solution towards the original classification (January) without preventing it from adjusting that classification in February as the data determines.

As mentioned above, in item 108, the documents in T2 are classified using the classifier 218. In item 114, the mean of these classified T2 documents in

document space D2 is used to create seeds in the D2 space using the centroid seed generator 220. New clusters are then generated in item 118 by the K-means cluster generator 222 to generate T2-D2 clusters 228.

Item 120 is an optional process used to add additional J clusters to the set of K seeds, depending on whether the user suspects additional concepts have been introduced in the new dataset T2, as determined by the comparator 230 (which compares the T1-D1 cluster and the T2-D2 cluster). If the new concepts found in item 120 are not useful (e.g. represent uninteresting subconcepts of the original K classes) then the result may be discarded or the additional new J-classes may be output 232.

At the users discretion, the additional J clusters may be added to capture any new, emerging concepts that may be contained in dataset T2. New J-class seeds may be chosen in any number of ways, including by random selection of data points from T2. In general, carefully choosing the new seeds to reflect areas of the data that are not well represented by the centroids generated in item 17 will produce the best results. The number of J additional clusters is a user determined parameter (just as K is a user determined parameter). The user determines the correct value through a process of trial and error. For example, if the user choose a value J=5 and saw a result where two of the additional cluster were spurious (i.e. a cluster is spurious if it contains no new concepts or too few examples to be meaningful) then the correct value of J would be 3. If all 5 new clusters contained

valuable concepts then the user might next try $J=10$, and so on until J was set large enough to create some spurious clusters.

The invention can be implemented, for example, as a computer program, written in the Java programming language and running on top of the Java virtual machine. An implementation which uses random sampling to throw away the poorest seeds is described below.

```
package com.ibm.cv.text;

import java.io.*;
import java.util.*;
import java.awt.*;
import com.ibm.cv.*;

public class AddClusters {

    public int numToAdd = 0;
    public int originalSize = 0;
    public KMeans k;
    public static int numAdditional = 10;

    //Given a KMeans clustering this method will add the given
    //number of clusters to it.

    public AddClusters(KMeans x, int added) {
        numToAdd = added;
        k = x;
        originalSize = k.nclusters;
        int numseeds = k.nclusters + numToAdd;
        k.nclusters = numseeds;
        seeds = new float[numseeds] [];

        for (int i=0; i<originalSize; i++) // add the original centroids to the new
            centroid list seeds[i] = k.centroids[i];

        for (int i=0; i<(numToAdd+numAdditional); i++) { //create a set of potential
            centroids int e = (int)(k.ndata*Math.random());
```

```

        float pseed[] =k.getData(e);
        potentialSeeds.addElement(pseed);
    }

5    getGoodSeeds(); // eliminate bad centroids and iteratively and put the remainder
        // in the seeds array
    k.centroids = seeds; // reset the centroids of the KMeans object
    k.run(); // run KMeans

10 }

    public void getGoodSeeds() {
        // iteratively remove the worst candidate seed until only the best candidates
        // remain.
15    for (int i=0; i<numAdditional; i++)
        replaceWorstSeed();

        for (int i=originalSize; i<k.nclusters; i++)
            seeds[i] = (float[])potentialSeeds.elementAt(i-originalSize);
20 }

    public void replaceWorstSeed() {
        int pos = findWorstSeed();
        potentialSeeds.removeElementAt(pos);
25 }

    public int findWorstSeed() {
        // returns the seed with the lowest data membership count.
        int counts[] = new int[potentialSeeds.size()];
30

        for (int i=0; i<k.ndata; i++) {//count the number of data elements going
            // to each seed
            int s = findNearestSeed(i)-originalSize;
            if (s>=0) counts[s]++; //ignore data elements that go to the original
35            // centroids
        }

        int result = Util.min(counts);
        return(Util.findPosition(result, counts));
40

    }

```



```

public int findNearestSeed(int d) {
    int best = -1;
    float bestval = 1.0F;
    // first check the new candidate centroids
5    for (int i=0; i<potentialSeeds.size0; i++) {
        float pseed[] = (float[])potentialSeeds.elementAt(i); // pseed =
        PotentialCentroid
        float val = k.getDistance(d,pseed); // calculate distance between point
        d and pseed
10    if (val<bestval) {
        best = i+originalSize;
        bestval = val;
    }
}
15 // then check the old centroids
for (int i=0; i<originalSize; i++) {
    float val = k.getDistance(d,k.centroids[i]); //calculate distance between d
    and original centroid
    if (val<bestval) {
20        best = i;
        bestval = val;
    }
}
// return the best centroid.
25 return(best);
}

```

While the overall methodology of the invention is described above, the invention can be embodied in any number of different types of systems and executed in any number of different ways, as would be known by one ordinarily skilled in the art. For example, as illustrated in Figure 3, a typical hardware configuration of an information handling/computer system in accordance with the invention preferably has at least one processor or central processing unit (CPU)

300. For example, the central processing unit 300 could include various

35 image/texture processing units, mapping units, weighting units, classification

units, clustering units, filters, adders, subtractors, comparators, etc. Alternatively, as would be known by one ordinarily skilled in the art given this disclosure, multiple specialized CPU's (or other similar individual functional units) could perform the same processing, mapping, weighting, classifying, clustering, filtering, adding, subtracting, comparing, etc.

The CPU 300 is interconnected via a system bus 301 to a random access memory (RAM) 302, read-only memory (ROM) 303, input/output (I/O) adapter 304 (for connecting peripheral devices such as disk units 305 and tape drives 306 to the bus 301), communication adapter 307 (for connecting an information handling system to a data processing network) user interface adapter 308 (for connecting peripherals 309-310 such as a keyboard, mouse, imager, microphone, speaker and/or other interface device to the bus 301), a printer 311, and display adapter 312 (for connecting the bus 301 to a display device 313). The invention could be implemented using the structure shown in Figure 3 by including the inventive method, described above, within a computer program stored on the storage device 305. Such a computer program would act on an image supplied through the interface units 309-310 or through the network connection 307. The system would then automatically segment the textures and output the same on the display 313, through the printer 311 or back to the network 307.

CLAIMS

What is claimed is:

- 5 1. A method of clustering documents in datasets comprising:
 clustering first documents and a first dataset to produce first document
 classes;
 creating centroid seeds based on said first document classes; and
 clustering second documents in a second dataset using said centroid seeds.
- 10 2. The method in claim 1, wherein said first dataset and said second dataset
 are related.
3. The method in claim 1, wherein said clustering of said first documents in
15 said first dataset comprises:
 forming a first dictionary of most common words in said first dataset;
 generating a first vector space model by counting, for each word in said
 first dictionary, a number of said first documents in which said word occurs; and
 clustering said first documents in said first dataset based on said first
20 vector space model.

4. The method in claim 3, further comprising generating a second vector space model by counting, for each word in said first dictionary, a number of said second documents in which said word occurs.

5 5. The method in claim 4, wherein said creating of said centroid seeds comprises:

classifying said second vector space model using said first document classes to produce a classified second vector space model; and

determining a mean of vectors in each class in said classified second
10 vector space model, wherein said mean comprises said centroid seeds.

6. The method in claim 4, further comprising:

forming a second dictionary of most common words in said second dataset;

15 generating a third vector space model by counting, for each word in said second dictionary, a number of said second documents in which said word occurs; and

clustering said documents in said second dataset based on said third vector space model to produce a second dataset cluster.

20

7. The method in claim 6, wherein said clustering of said second documents in said second dataset using said centroid seeds produces an adapted dataset cluster and said method further comprises:

comparing classes in said adapted dataset cluster to classes in said second
5 dataset cluster; and
adding classes to said adapted dataset cluster based on said comparing.

8. A system for clustering documents in datasets comprising:

a storage having a first dataset and a second dataset;
10 a cluster generator operative to cluster first documents in said first dataset and produce first document classes; and
a centroid seed generator operative to generate centroid seeds based on said first document classes,
wherein said cluster generator clusters second documents in said second
15 dataset using said centroid seeds.

9. The system in claim 8, wherein said first dataset and said second dataset are related.

10. The system in claim 8, further comprising:

a dictionary generator adapted to generate a first dictionary of most common words in said first dataset; and

a vector space model generator adapted to generate a first vector space model by counting, for each word in said first dictionary, a number of said first documents in which said word occurs,

wherein said cluster generator clusters said documents in said first dataset based on said first vector space model.

11. The system in claim 10, wherein said vector space model generator generates a second vector space model by counting, for each word in said first dictionary, a number of said second documents in which said word occurs.

12. The system in claim 11, further comprising a classifier adapted to classify said second documents in said second vector space model using said first document classes to produce a classified second vector space model and adapted to determine a mean of vectors in each class in said classified second vector space model, wherein said mean comprises said centroid seeds.

13. The system in claim 11, wherein:

said dictionary generator is adapted to generate a second dictionary of most common words in said second dataset,

said vector space model generator is adapted to generate a third vector space model by counting, for each word in said second dictionary, a number of said second documents in which said word occurs, and

said cluster generator is adapted to cluster said second documents in said second dataset based on said third vector space model to produce a second dataset cluster.

10

14. The system in claim 13, wherein said cluster generator is adapted to produce an adapted dataset cluster by clustering said second documents in said second dataset using said centroid seeds and said system further comprises:

a comparator adapted to compare classes in said adapted dataset cluster to classes in said second dataset cluster and add classes to said adapted dataset cluster based on said comparing.

15. A method of clustering documents in a first dataset having first documents and a related second dataset having second documents, said method comprising:

clustering said first documents to produce first document classes;

generating a vector space model of said second documents;

classifying said vector space model of said second documents using said first document classes to produce a classified vector space model; and
determining a mean of vectors in each class in said classified vector space model to produce centroid seeds; and
5 clustering said second documents using said centroid seeds.

16. The method in claim 15, wherein said vector space model comprises a second vector space model and said clustering of said first documents in said first data comprises:

10 forming a first dictionary of most common words in said first dataset; and
generating a first vector space model by counting, for each word in said first dictionary, a number of said first documents in which said word occurs,
wherein said clustering of said first documents in said first dataset is based on said first vector space model.

15 17. The method in claim 16, wherein said generating of said second vector space model comprises counting, for each word in said first dictionary, a number of said second documents in which said word occurs.

18. The method in claim 17, further comprising:

forming a second dictionary of most common words in said second dataset;

generating a third vector space model by counting, for each word in said
5 second dictionary, a number of said second documents in which said word occurs;
and

clustering said documents in said second dataset based on said third vector
space model to produce a second dataset cluster.

10 19. The method in claim 18, wherein said clustering of said second documents
in said second dataset using said centroid seeds produces an adapted dataset
cluster and said method further comprises:

comparing classes in said adapted dataset cluster to classes in said second
dataset cluster; and

15 adding classes to said adapted dataset cluster based on said comparing.

20. A method of clustering documents in related datasets comprising:

forming a first dictionary of most common words in a first dataset;

generating a first vector space model by counting, for each word in said
20 first dictionary, a number of said first documents in which said word occurs; and

clustering said first documents in said first dataset based on said first
vector space model to produce first document classes;

generating a second vector space model by counting, for each word in said
first dictionary, a number of said second documents in which said word occurs;

5 classifying said second documents in said second vector space model
using said first document classes to produce a classified second vector space
model;

determining a mean of vectors in each class in said classified second
vector space model to produce centroid seeds; and

10 clustering second documents in a second dataset using said centroid seeds

21. The method in claim 20, further comprising:

forming a second dictionary of most common words in said second
dataset;

15 generating a third vector space model by counting, for each word, in said
second dictionary, a number of said second documents in which said word occurs;
and

clustering said documents in said second dataset based on said third vector
space model to produce a second dataset cluster.

20

22. The method in claim 21, wherein said clustering of said second documents in said second dataset using said centroid seeds produces an adapted dataset cluster and said method further comprises:

5 comparing classes in said adapted dataset cluster to classes in said second dataset cluster; and
adding classes to said adapted dataset cluster based on said comparing.

23. A program device readable by machines, tangibly embodying a program of instructions executable by the machine to perform said method of clustering documents in datasets comprising:

10 clustering first documents and a first dataset to produce first document classes;
creating centroid seeds based on said first document classes; and
clustering second documents in a second dataset using said centroid seeds.

15 24. A program device readable by machines, tangibly embodying a program of instructions executable by the machine to perform said method in claim 23, wherein said first dataset and second dataset are related.

25. A program device readable by machines, tangibly embodying a program of instructions executable by the machine to perform said method in claim 23, wherein said clustering of said first documents in said first dataset comprises:

forming a first dictionary of most common words in said first dataset;

5 generating a first vector space model by counting, for each word in said first dictionary, a number of said first documents in which said word occurs; and clustering said first documents in said first dataset based on said first vector space model.

10 26. A program device readable by machines, tangibly embodying a program of instructions executable by the machine to perform said method in claim 25, further comprising generating a second vector space model by counting, for each word in said first dictionary, a number of said second documents in which said word occurs.

15 27. A program device readable by machines, tangibly embodying a program of instructions executable by the machine to perform said method in claim 26, wherein said creating of said centroid seeds comprises:

classifying said second vector space model using said first document

20 classes to produce a classified second vector space model; and

determining a mean of vectors in each class in said classified second vector space model, wherein said mean comprises said centroid seeds.

28. A program device readable by machines, tangibly embodying a program of instructions executable by the machine to perform said method in claim 26, further comprising:

forming a second dictionary of most common words in said second dataset;

generating a third vector space model by counting, for each word in said second dictionary, a number of said second documents in which said word occurs; and

clustering said documents in said second dataset based on said third vector space model to produce a second dataset cluster.

29. A program device readable by machines, tangibly embodying a program of instructions executable by the machine to perform said method in claim 28, wherein said clustering of said second documents in said second dataset using said centroid seeds produces an adapted dataset cluster and said method further comprises:

comparing classes in said adapted dataset cluster to classes in said second dataset cluster; and

adding classes to said adapted dataset cluster based on said comparing.

5 30. A program device readable by machines, tangibly embodying a program of instructions executable by the machine to perform a system for clustering documents in datasets comprising:

a storage having a first dataset and a second dataset;

a cluster generator operative to cluster first documents in said first dataset

10 and produce first document classes; and

a centroid seed generator operative to generate centroid seeds based on said first document classes,

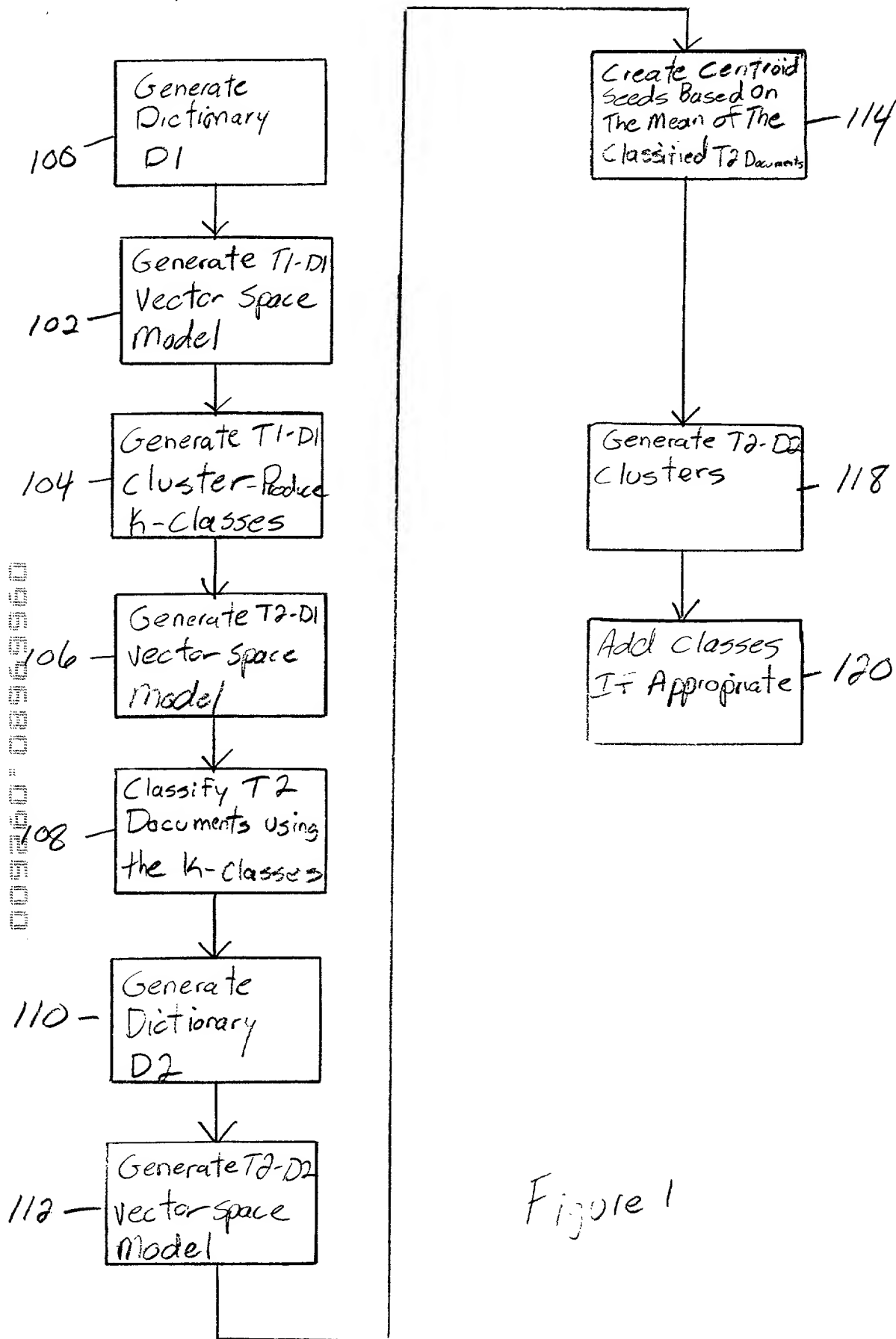
wherein said cluster generator clusters second documents in said second dataset using said centroid seeds.

15

METHOD FOR ADAPTING A K-MEANS TEXT CLUSTERING TO EMERGING DATA

ABSTRACT

A method and structure for clustering documents in datasets which include
5 clustering first documents and a first dataset to produce first document classes,
creating centroid seeds based on the first document classes, and clustering second
documents in a second dataset using the centroid seeds, wherein the first dataset
and the second dataset are related. The clustering of the first documents in the
first dataset forms a first dictionary of most common words in the first dataset and
10 generates a first vector space model by counting, for each word in the first
dictionary, a number of the first documents in which the word occurs, and clusters
the first documents in the first dataset based on the first vector space model, and
further generates a second vector space model by counting, for each word in the
first dictionary, a number of the second documents in which the word occurs.
15 Creation of the centroid seeds includes classifying second vector space model
using the first document classes to produce a classified second vector space model
and determining a mean of vectors in each class in the classified second vector
space model, the mean includes the centroid seeds.



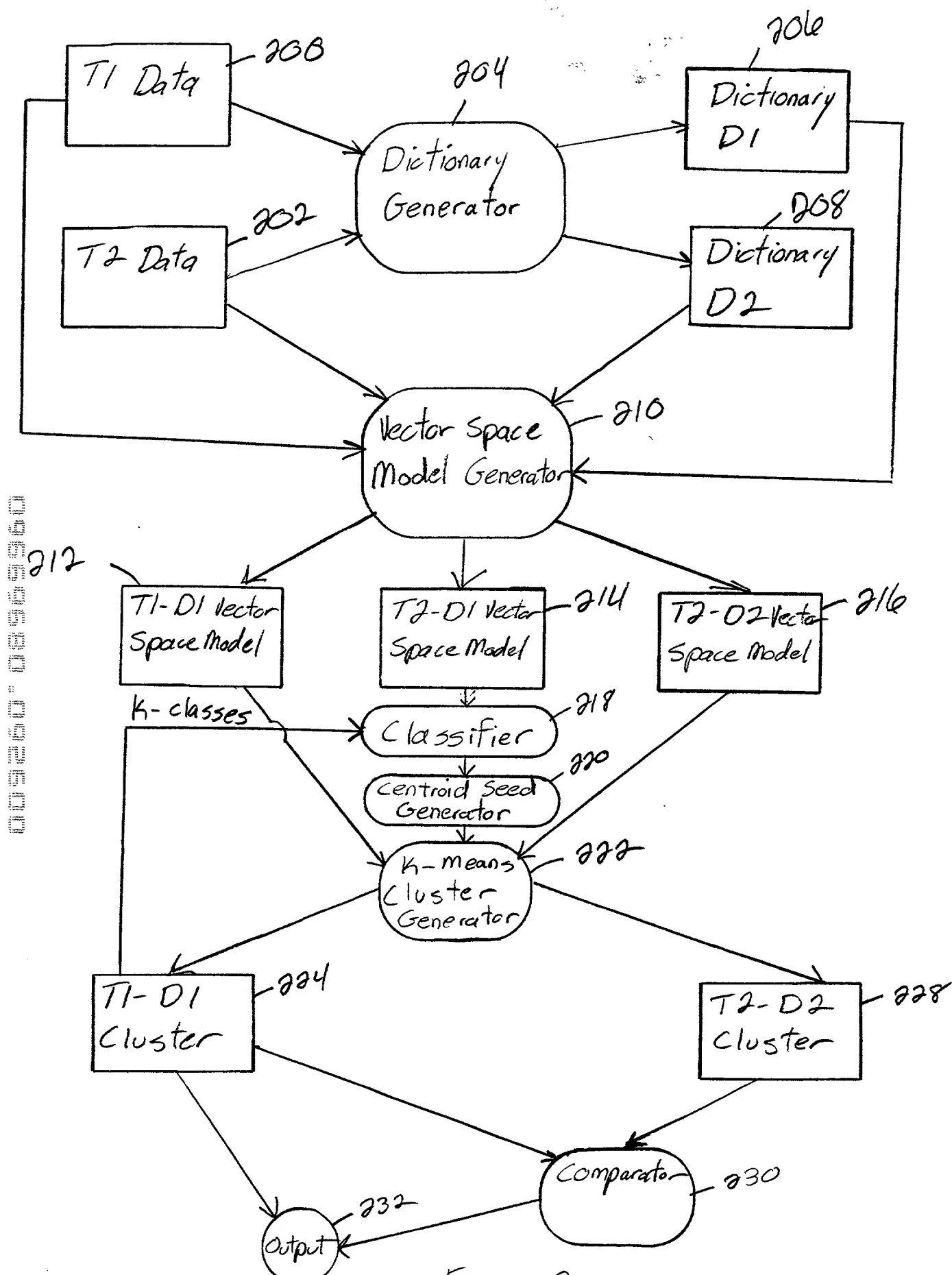


Figure 2

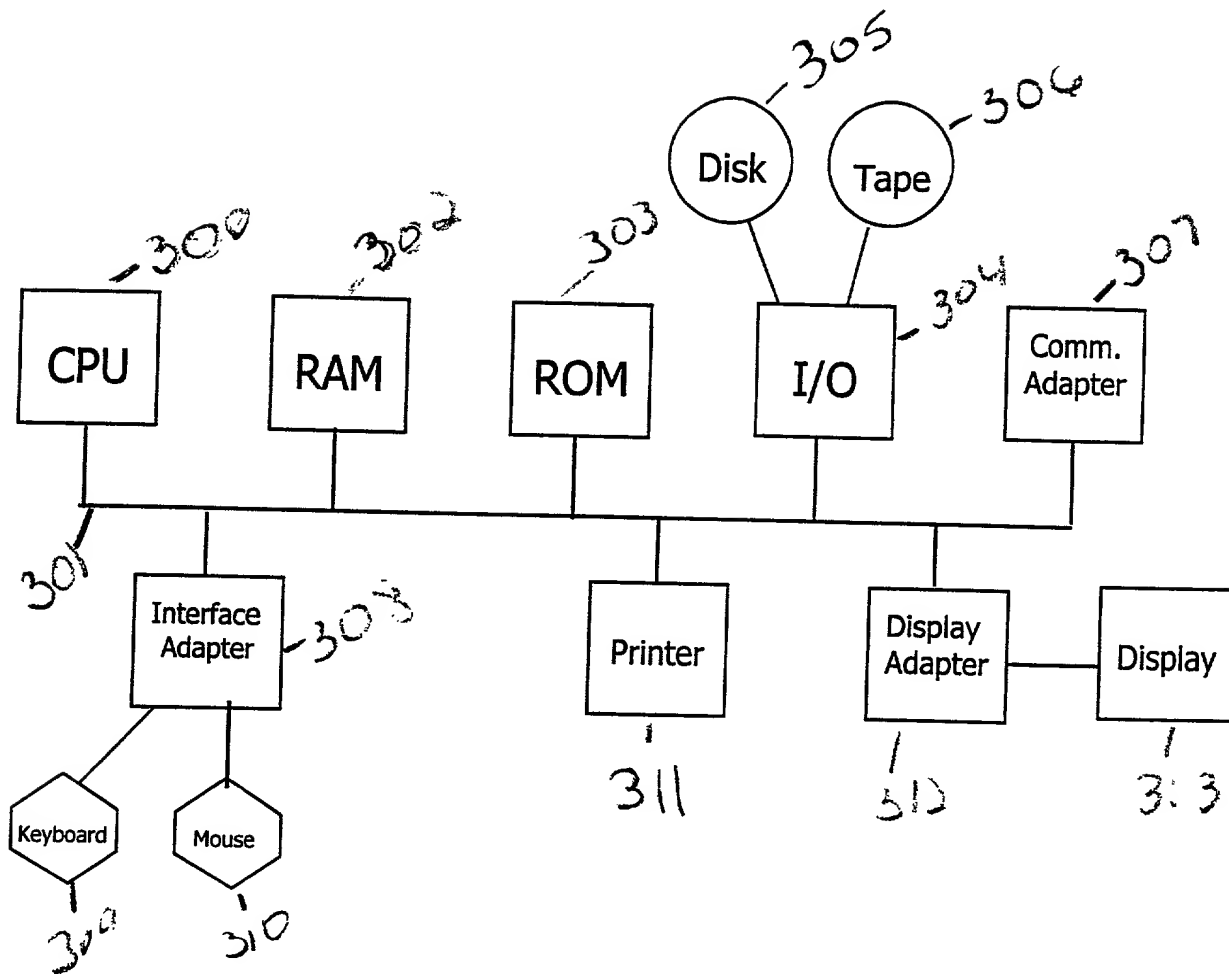


Figure 3

DECLARATION AND POWER OF ATTORNEY

As a below named inventor, I hereby declare that:

My residence, post office address and citizenship are as stated below next to my name; I believe I am the original, first and sole inventor (if only one name is listed below) or an original, first and joint inventor (if plural names are listed below) of the subject matter which is claimed and for which a patent is sought on the invention entitled: A METHOD FOR ADAPTING A K-MEANS TEXT CLUSTERING TO EMERGING DATA

the specification of which:
(check one)

☒ is attached hereto.

☐ was filed on _____, as Application Serial No. _____ and was amended on _____.

I hereby state that I have reviewed and understand the contents of the above identified specification, including the claims, as amended by any amendment referred to above.

I acknowledge the duty to disclose information which is material to the patentability of this application in accordance with Title 37, Code of Federal Regulations, § 1.56.

I hereby claim foreign priority benefits under Title 35, United States Code, § 119 of any foreign application(s) for patent or inventor's certificate listed below and have also identified below any foreign application for patent or inventor's certificate having a filing date before that of the application on which priority is claimed:

Prior Foreign Application(s):

Number	Country	Day/Month/Year	Priority Claimed
--------	---------	----------------	------------------

I hereby claim the benefit under Title 35, United States Code, § 120 of any United States application(s) listed below and, insofar as the subject matter of each of the claims of this application is not disclosed in the prior United States application in the manner provided by the first paragraph of Title 35, United States Code, § 112, I acknowledge the duty to disclose material information as defined in Title 37, Code of Federal Regulations, § 1.56 which occurred between the filing date of the prior application and the national or PCT international filing date of this application:

Prior U.S. Applications:

Serial No.	Filing Date	Status
------------	-------------	--------

I hereby declare that all statements made herein of my own knowledge are true and that all statements made on information and belief are believed to be true; and further that these statements were made with the knowledge that willful false statements and the like so made are punishable by fine or imprisonment, or both, under Section 1001 of Title 18 of the United States Code and that such willful false statements may jeopardize the validity of the application or any patent issued thereon.

As a named inventor, I hereby appoint the following attorneys and/or agents to prosecute this application and transact all business in the Patent and Trademark Office connected therewith: We hereby appoint Khanh Q. Tran, Registration No. 41,352, Thomas R. Berthold, Registration No. 28,689, Marc McSwain, Registration No. 44,929, Alison D. Mortinger, Registration No. 39,306, Frederick W. Gibb, III, Registration No. 37,629, and Sean M. McGinn, Registration No. 34,386 to prosecute this application and transact all business in the United States Patent and Trademark Office connected therewith.

Send all correspondence to: McGinn & Gibb, P.C., 1701 Clarendon Boulevard, Suite 100, Arlington, Virginia 22209. Customer No. 21254

Telephone calls should be directed to Frederick W. Gibb, III, McGinn & Gibb, P.C. at (703) 294-6699.

(1) Inventor: William Scott Spangler

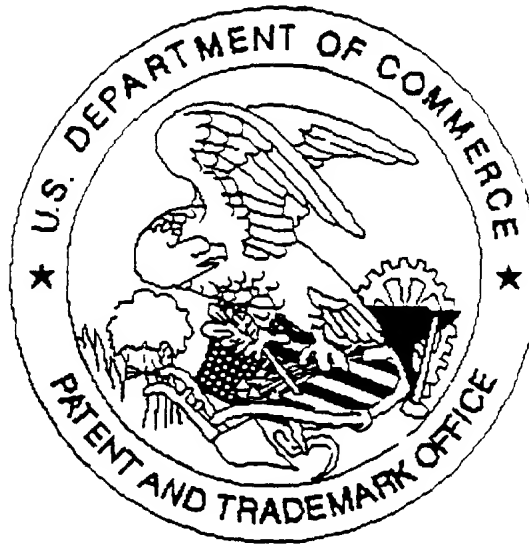
Signature: William Scott Spangler Date: 9/11/00

Residence: 12840 Stevens Court, San Martin, CA 95046

Citizenship: USA

Post Office Address: Same as Residence

United States Patent & Trademark Office
Office of Initial Patent Examination -- Scanning Division



SCANNED, # 24

Application deficiencies were found during scanning:

☐ Page(s) Bof 3 of Trans with were not present
for scanning. (Document title)

☐ Page(s) _____ of _____ were not present
for scanning. (Document title)

☐ Scanned copy is best available.